

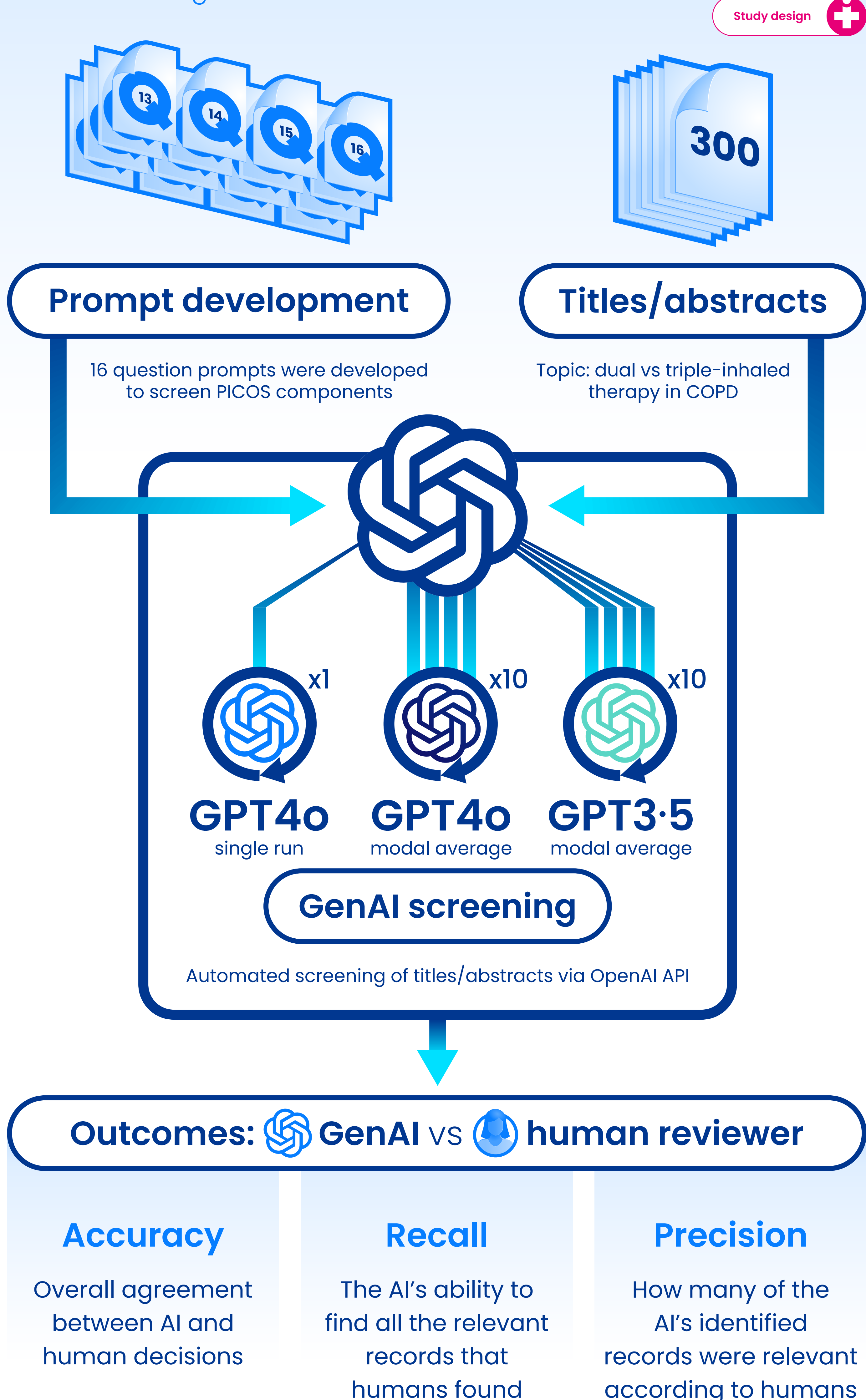
Automating screening using generative AI (GenAI) could speed up literature reviews, but performance may depend on model type, topic complexity, and prompt design

Background

- Literature screening for systematic literature reviews (SLRs) is time-consuming
- GenAI may streamline this process, but there are limited data on its performance
- This study evaluated the level of agreement between genAI and a human reviewer when screening 300 titles and abstracts from a previous SLR (Fig 1)

Method

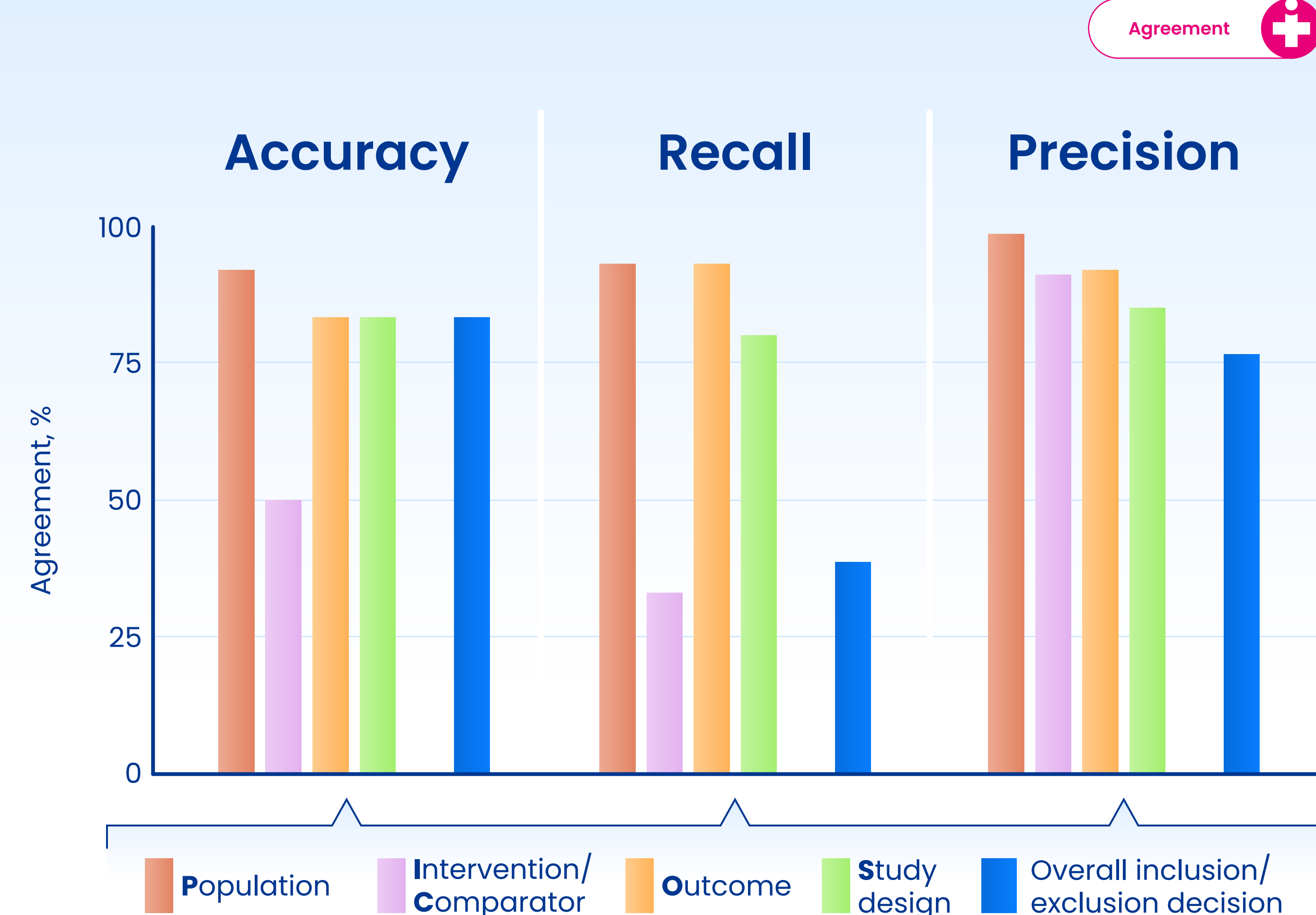
Fig 1. Approaches to assess the performance of genAI vs human screening of literature



Results

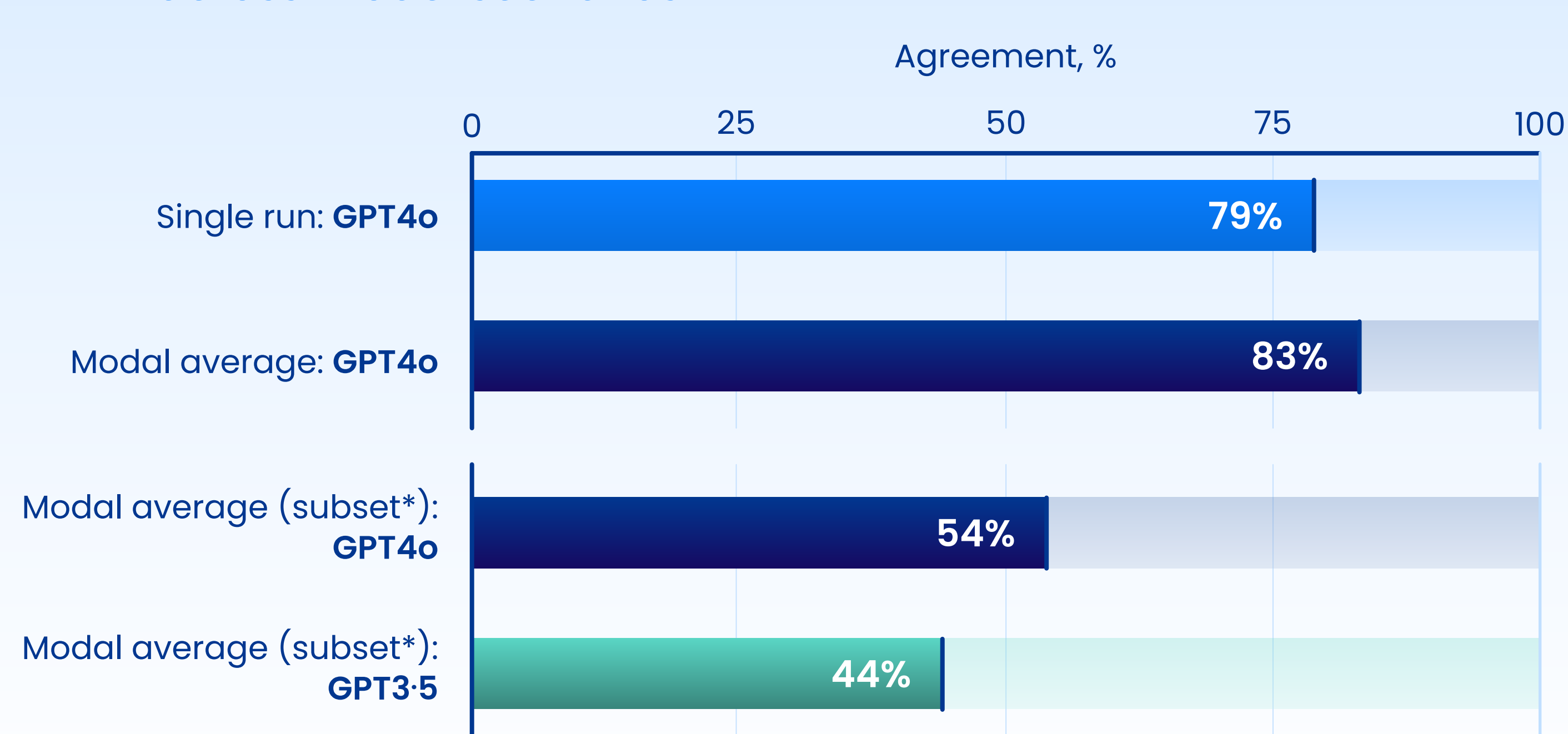
- GenAI showed high agreement with the human reviewer on Population, Outcome and Study design, but low agreement on the Intervention/Comparator domain (Fig 2)
- Accuracy and precision were relatively high for the overall inclusion/exclusion decision, but recall was low

Fig 2. Agreement between single-run GPT4o and human reviewer across PICOS domains



- Compared with the single-run approach, taking the modal answer of 10 runs with GPT4o resulted in slightly higher accuracy for the overall inclusion/exclusion decision (Fig 3)
- Conversely, accuracy was lower for the GPT3.5 subset than with the GPT4o subset

Fig 3. Overview of accuracy for overall inclusion/exclusion decision across model scenarios



*In a selected subset of 50/300 records

- Overall, agreement between genAI and the human reviewer (accuracy) was similar to what would be expected between two humans!